

Efficient Fault-Tolerant Cluster-Sending: *Reliable and Efficient Communication between Byzantine Fault-Tolerant Clusters*

*Jelle Hellings*¹ Mohammad Sadoghi²

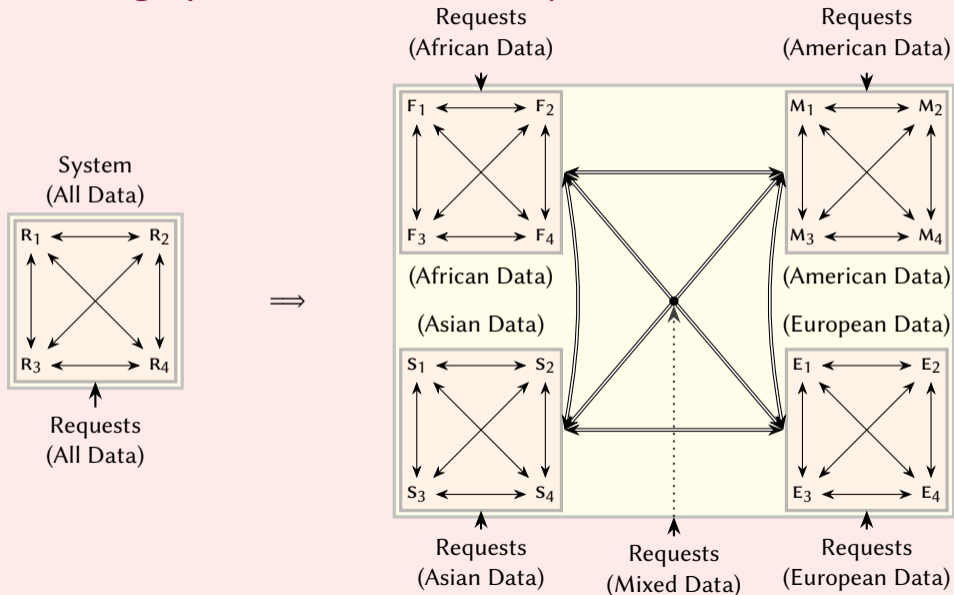
¹Department of Computing and Software
McMaster University



²Exploratory Systems Lab
Department of Computer Science
University of California, Davis



Motivation: High-performance resilient system



Blockchain interoperability in the wild

Permissionless blockchains

- ▶ Bitcoin, Ethereum,
- ▶ Proof-of-Work, Proof-of-Stake,
- ▶ BTC Relay.
- ▶ Cosmos, Polkadot.
- ▶ Atomic swaps, atomic commitment, cross-chain deals,

Blockchain interoperability in the wild

Permissionless blockchains

- ▶ Bitcoin, Ethereum,
- ▶ Proof-of-Work, Proof-of-Stake,
- ▶ BTC Relay.
- ▶ Cosmos, Polkadot.
- ▶ Atomic swaps, atomic commitment, cross-chain deals,

Permissioned blockchains

- ▶ ResilientDB, BFT-SMaRt,
- ▶ Practical Byzantine Fault Tolerance.
- ▶ Steward, GeoBFT.
- ▶ AHL, ByShard, Chainspace, Cerberus, RingBFT.
- ▶ *The cluster-sending problem.*

The Byzantine cluster-sending problem

The problem of sending a value v from a cluster C_1 to a cluster C_2 such that

- ▶ all non-faulty replicas in C_2 *RECEIVE* the value v ;
- ▶ all non-faulty replicas in C_1 *CONFIRM* that the value v was received; and
- ▶ C_2 only receives a value v if all non-faulty replicas in C_1 *AGREE* upon sending v .

Requirements to provide reliable communication between clusters with Byzantine replicas.

Global communication versus local communication

Straightforward cluster-sending solution (crash failures)

Pair-wise broadcasting with $(f_1 + 1)(f_2 + 1) \approx f_1 \times f_2$ messages.

	<i>Ping round-trip times (ms)</i>						<i>Bandwidth (Mbit/s)</i>					
	OR	IA	Mont.	BE	TW	Syd.	OR	IA	Mont.	BE	TW	Syd.
Oregon	≤ 1	38	65	136	118	161	7998	669	371	194	188	136
Iowa		≤ 1	33	98	153	172		10004	752	243	144	120
Montreal			≤ 1	82	186	202			7977	283	111	102
Belgium				≤ 1	252	270				9728	79	66
Taiwan					≤ 1	137					7998	160
Sydney						≤ 1						7977

Lower bounds for cluster-sending: Example

$$\mathbf{n}_1 = 15$$

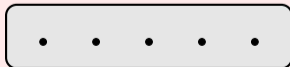
$$\mathbf{f}_1 = 7$$

$$\mathbf{n}_2 = 5$$

$$\mathbf{f}_2 = 2$$

Claim (crash failures)

Any correct algorithm needs to send at least 14 messages.



Lower bounds for cluster-sending: Example

$$n_1 = 15$$

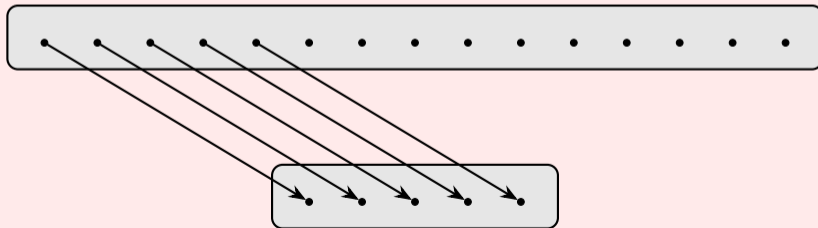
$$f_1 = 7$$

$$n_2 = 5$$

$$f_2 = 2$$

Claim (crash failures)

Any correct algorithm needs to send at least 14 messages.



Lower bounds for cluster-sending: Example

$$n_1 = 15$$

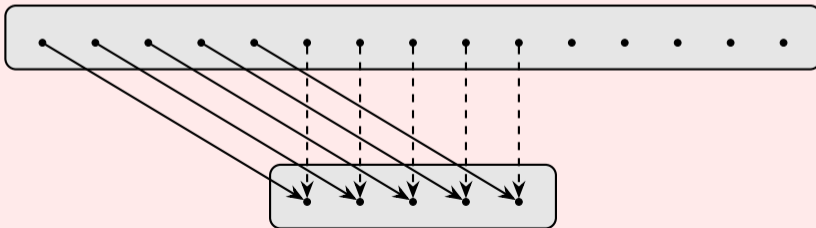
$$n_2 = 5$$

$$f_1 = 7$$

$$f_2 = 2$$

Claim (crash failures)

Any correct algorithm needs to send at least 14 messages.



Lower bounds for cluster-sending: Example

$$n_1 = 15$$

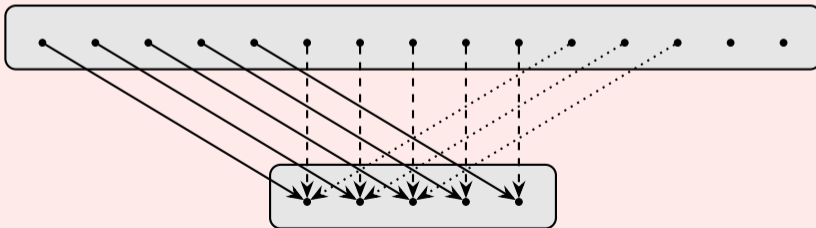
$$f_1 = 7$$

$$n_2 = 5$$

$$f_2 = 2$$

Claim (crash failures)

Any correct algorithm needs to send at least 14 messages.



Lower bounds for cluster-sending: Example

$$n_1 = 15$$

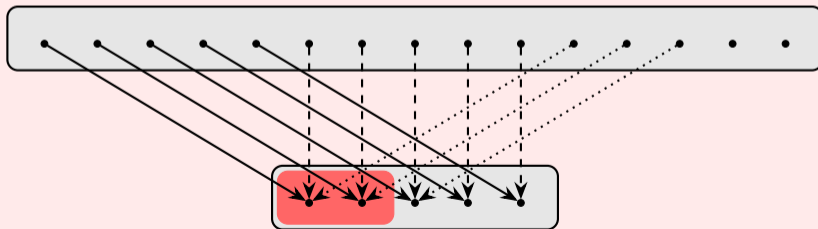
$$n_2 = 5$$

$$f_1 = 7$$

$$f_2 = 2$$

Claim (crash failures)

Any correct algorithm needs to send at least 14 messages.



Lower bounds for cluster-sending: Example

$$n_1 = 15$$

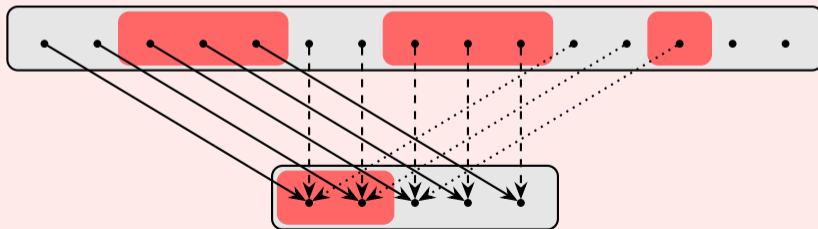
$$f_1 = 7$$

$$n_2 = 5$$

$$f_2 = 2$$

Claim (crash failures)

Any correct algorithm needs to send at least 14 messages.



Lower bounds for cluster-sending: Results

Theorem (Cluster-sending lower bound, simplified)

We need to exchange $\max(\mathbf{n}_1, \mathbf{n}_2)$ messages to do cluster-sending.

Theorem (Cluster-sending lower bound, crash failures)

Assume $\mathbf{n}_1 \geq \mathbf{n}_2$ and let

$$q = (\mathbf{f}_1 + 1) \operatorname{div} \mathbf{n}_2;$$

$$r = (\mathbf{f}_1 + 1) \operatorname{mod} \mathbf{n}_2.$$

We need to exchange at least $q\mathbf{n}_2 + r + \mathbf{f}_2 \operatorname{sgn} r \approx \mathbf{n}_1$ messages to do cluster-sending.

Cluster-sending via bijective sending

Protocol for the sending cluster C_1 :

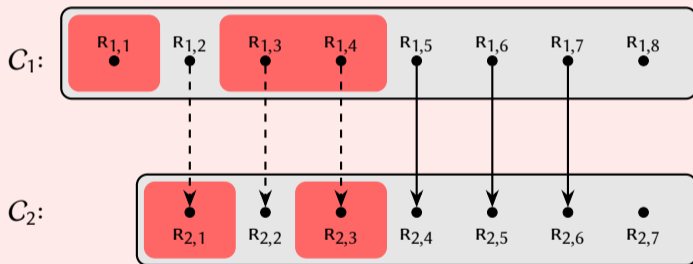
- 1: All replicas in $\text{nf}(C_1)$ agree on v and construct $\langle v \rangle_{C_1}$.
- 2: Choose replicas $S_1 \subseteq C_1$ with $\mathbf{n}_{S_1} = \mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$.
- 3: Choose replicas $S_2 \subseteq C_2$ with $\mathbf{n}_{S_2} = \mathbf{f}_{C_1} + \mathbf{f}_{C_2} + 1$.
- 4: Choose a bijection $b : S_1 \rightarrow S_2$.
- 5: **for** $R_1 \in S_1$ **do**
- 6: R_1 sends $(v, \langle v \rangle_{C_1})$ to $b(R_1)$.

Protocol for the receiving cluster C_2 :

- 7: **event** $R_2 \in \text{nf}(C_2)$ receives $(w, \langle w \rangle_{C_1})$ from $R_1 \in C_1$ **do**
 - 8: Broadcast $(w, \langle w \rangle_{C_1})$ to all replicas in C_2 .
 - 9: **event** $R'_2 \in \text{nf}(C_2)$ receives $(w, \langle w \rangle_{C_1})$ from $R_2 \in C_2$ **do**
 - 10: R'_2 considers w *received*.
-

Cluster-sending via bijective sending—visualized

$$\mathbf{n}_1 = 8, \mathbf{n}_2 = 7, \mathbf{f}_1 = 3, \mathbf{f}_2 = 2, \sigma = 6$$



Cluster-sending: Can we do better?

Pessimistic

No: these algorithms are worst-case optimal.

Cannot do better than *linear communication* in the size of the clusters.

Cluster-sending: Can we do better?

Pessimistic

No: these algorithms are worst-case optimal.

Cannot do better than *linear communication* in the size of the clusters.

Probabilistic approaches

Yes: if we randomly choose sender and receiver, then we often do much better!

Cluster-sending via probabilistic cluster-sending

Protocol for the sending cluster C_1 :

- 1: All replicas in $\text{nf}(C_1)$ agree on v and construct $\langle v \rangle_{C_1}$.
 - 2: **repeat**
 - 3: Choose replicas $(R_1, R_2) \in C_1 \times C_2$, fully at random.
 - 4: CS-STEP(R_1, R_2, v).
 - 5: *Wait* until the CS-STEP should finish.
 - 6: **until** C_1 reaches consensus on $\langle \text{proof}(\langle \text{send}(v) \rangle_{C_2}) \rangle_{C_1}$.
-

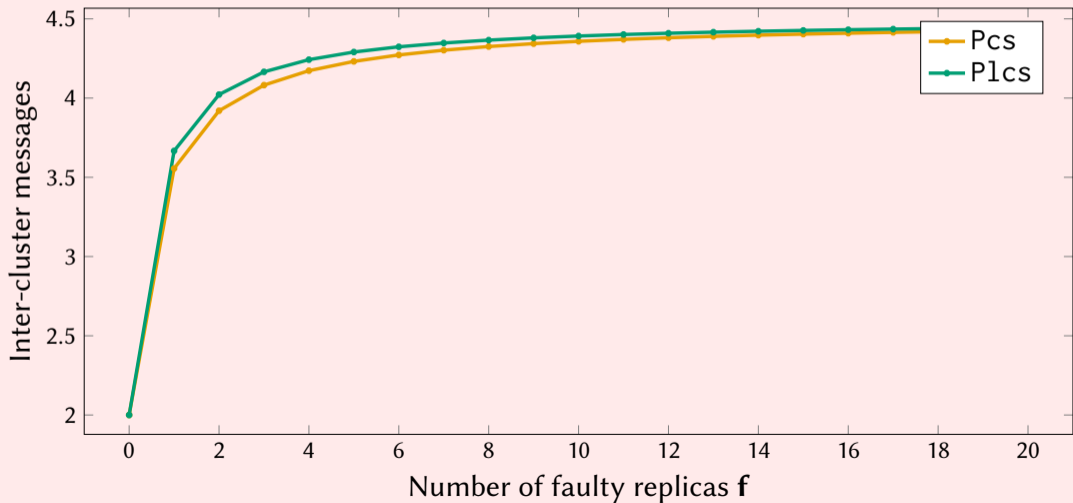
Cluster-sending via probabilistic cluster-sending

Protocol for the sending cluster C_1 :

- 1: All replicas in $\text{nf}(C_1)$ agree on v and construct $\langle v \rangle_{C_1}$.
 - 2: Let $(S_1, S_2) := \Phi(C_1, C_2)$.
 - 3: Choose $(P_1, P_2) \in \text{perms}(S_1) \times \text{perms}(S_2)$ fully at random.
 - 4: $i := 0$.
 - 5: **repeat**
 - 6: CS-STEP($P_1[i], P_2[i], v$).
 - 7: *Wait* until the CS-STEP should finish.
 - 8: $i := i + 1$.
 - 9: **until** C_1 reaches consensus on $\langle \text{proof}(\langle \text{send}(v) \rangle_{C_2}) \rangle_{C_1}$.
-

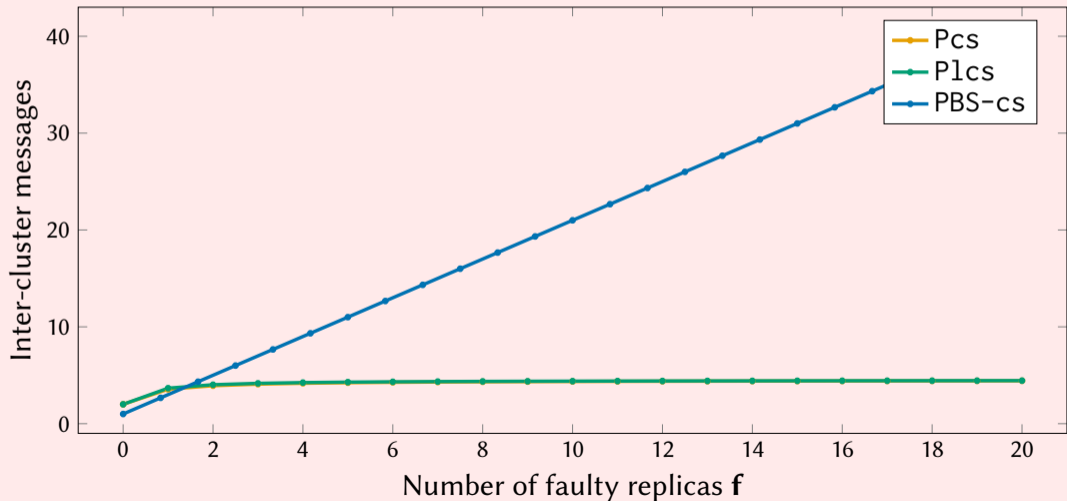
Performance of cluster-sending solutions

($n = 3f + 1$ replicas per cluster)



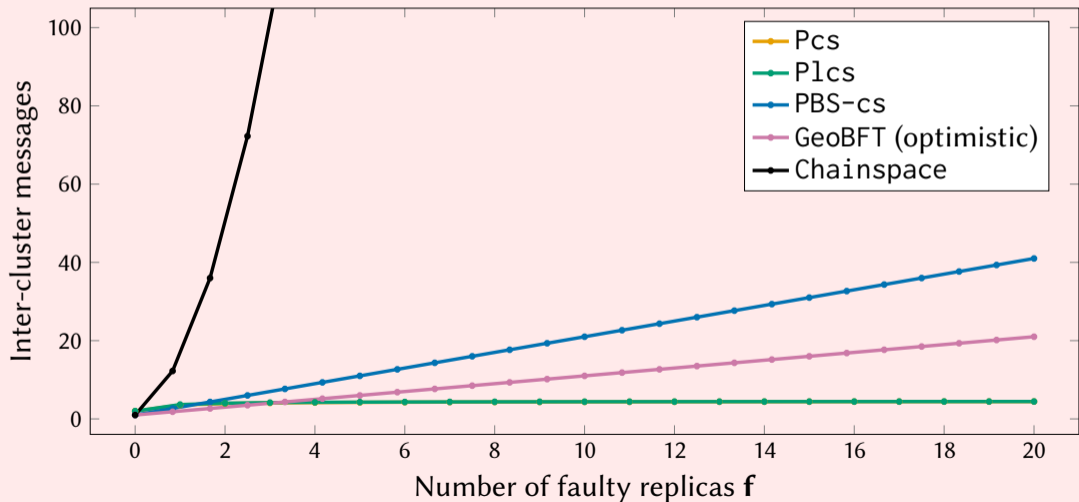
Performance of cluster-sending solutions

($n = 3f + 1$ replicas per cluster)



Performance of cluster-sending solutions

($n = 3f + 1$ replicas per cluster)



Cluster-sending in practice

Wide-area consensus GeoBFT uses *optimistic* cluster-sending techniques.

Cluster-sending in practice

Wide-area consensus GeoBFT uses *optimistic* cluster-sending techniques.

Sharded resilient databases ByShard, Cerberus, and RingBFT.

Cluster-sending in practice

Wide-area consensus GeoBFT uses *optimistic* cluster-sending techniques.

Sharded resilient databases ByShard, Cerberus, and RingBFT.
AHL and Chainspace use cluster-sending *informally*.

Conclusion

More information

<https://jhellings.nl>

Technical reports:

- ▶ <https://arxiv.org/abs/1908.01455>.
- ▶ <https://arxiv.org/abs/2108.08541>.